

A COMPARISON OF TWO APPROACHES: MAXIMUM ENTROPY ON THE MEAN (MEM) AND BAYESIAN ESTIMATION (BAYES) FOR INVERSE PROBLEMS

ALI MOHAMMAD-DJAFARI
Laboratoire des Signaux et Systèmes (CNRS-ESE-UPS)
École Supérieure d'Électricité,
Plateau de Moulon, 91192 Gif-sur-Yvette, France.
E-mail: djafari@lss.supelec.fr

Abstract. To handle with inverse problems, two probabilistic approaches have been proposed: the maximum entropy on the mean (MEM) and the Bayesian estimation (BAYES). The main object of this presentation is to compare these two approaches which are in fact two different inference procedures to define the solution of an inverse problem as the optimizer of a compound criterion.

Key words: Inverse problems, Maximum Entropy on the Mean, Bayesian inference, Convex analysis

1. Introduction

Inverse problems arises in many areas of science and engineering. In fact, rarely, we can measure directly a quantity x and, in general, the unobserved interested x is related to the measured quantity y via a model. In many area this model can be written in the general form $y = \mathcal{A}(x) + n$ or in the discrete case:

$$\mathbf{y} = \mathbf{A}(\mathbf{x}) + \mathbf{n}, \quad (1)$$

where \mathbf{y} stands for the data, \mathbf{x} for the unknown variables and \mathbf{n} for the errors (modeling and noise). Since Newton and Gauss, one tries to define a solution to this problem as the optimizer of a criterion, for example the Least Squares (LS):

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \{ \|\mathbf{y} - \mathbf{A}(\mathbf{x})\|^2 \}. \quad (2)$$

But the inverse problems are, in general, ill-posed and the LS criterion may not have a unique optimum or this solution may be very sensitive to noise. Since Tikhonov [?], the regularization theory became the main approach to give a satisfactory solution by defining it as the optimizer of a compound criterion:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \{ J(\mathbf{x}) \} \quad \text{with} \quad J(\mathbf{x}) = Q(\mathbf{x}) + \lambda \Omega(\mathbf{x}) = \|\mathbf{y} - \mathbf{A}(\mathbf{x})\|^2 + \lambda \|\mathbf{D}\mathbf{x}\|^2 \quad (3)$$

or in its more general forms [?]:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \{ J(\mathbf{x}) \} \quad \text{with} \quad J(\mathbf{x}) = Q(\mathbf{y} - \mathbf{A}(\mathbf{x})) + \lambda \Omega(\mathbf{x}, \mathbf{m}). \quad (4)$$

The questions then raised on how to choose the functionals Q and Ω and the regularization parameter λ and the default solution \mathbf{m} .

The probabilistic approaches started to give partial answers to this request. In particular in the Bayesian estimation approach and the maximum a posteriori (MAP) estimate:

$$\hat{\mathbf{x}} = \arg \max_{\mathbf{x}} \{p(\mathbf{x}|\mathbf{y})\} = \arg \min_{\mathbf{x}} \{-\log p(\mathbf{x}|\mathbf{y})\} = \arg \min_{\mathbf{x}} \{-\log p(\mathbf{y}|\mathbf{x}) - \log p(\mathbf{x})\}, \quad (5)$$

this choice is: $Q = -\log p(\mathbf{y}|\mathbf{x})$ and $\lambda\Omega(\mathbf{x}) = -\log p(\mathbf{x})$. This approach just pushed a little farther the questions which became how to translate our prior knowledge into a probability law and how to determine their parameters. Even, nowadays, there are many tools for the estimation of the hyperparameters [?], the main question on how to translate some knowledge about \mathbf{x} into a probability distribution stays without a complete answer. The maximum entropy (ME) principle gave partial answers [?, ?, ?]. See also [?] for an extensive discussed bibliography.

At the same time, many authors used the ME principle to find unique solutions to linear inverse problems by considering \mathbf{x} as a distribution and the data \mathbf{y} as linear constraints on them. Then, assuming that the data constraints are satisfied by a non empty set of solutions, a unique solution is chosen by maximizing the entropy:

$$-\sum_j x_j \log x_j \quad \text{or} \quad -\sum_j x_j \log \left[\frac{x_j}{m_j} - (x_j - m_j) \right], \quad (6)$$

where \mathbf{m} is default solution. See for example [?, ?] and the cited references. However, even if in these methods, thanks to convex analysis and Lagrangian techniques, the constrained optimization of 6 can be replaced by an equivalent unconstrained optimization, the obtained solutions satisfy the uniqueness condition of well-posedness but not the stability one [?, ?, ?].

Recently, some authors [?, ?, ?, ?] used the ME principle in a different way by considering \mathbf{x} not as a distribution but as the mean value of a random vector \mathbf{X} and the data as the constraints on its distribution $dP(\mathbf{x})$. Then, the ME principle is used to define it uniquely and finally the solution $\hat{\mathbf{x}}$ is defined as the expected value of this ME distribution.

Following these authors, some others used, commented and analyzed extensively these ideas [?, ?, ?, ?, ?, ?], [?, ?, ?, ?, ?] and [?, ?, ?, ?]. However, in all these works, the data \mathbf{y} were considered as exact constraints and the errors on the data were either neglected or partially taken account of. (See however new developments in [?].)

More recently, some authors who were more faced with real applications, [?, ?, ?, ?] followed the same idea, but by fixing themselves as the objective to use these ideas for describing the solution as the optimizer of a combined convex criteria such as (4) and more on a constructive way to determine these functionals.

The objective of this paper is to make a comparison of the Bayesian approach which we call hereafter BAYES and the maximum entropy in the mean which we refer to as MEM. This comparison is done very pragmatically and is based on the

understanding of the author who does not have pretension to know all the details of the both approaches and will be happy to discuss all the following discussions with the pro of the approaches.

2. Maximum entropy on the mean approach

2.1. BASICS

The main references to the basics of this approach are [?, ?]. The original idea and first applications in crystallography are given in [?, ?]. More details and extensions are given in [?, ?, ?, ?]. The mathematical aspects of convex analysis and duality theorems are given in [?, ?, ?, ?, ?, ?].

The following resumes the different steps of the approach:

- Consider a set \mathcal{C} , assume that $\mathbf{x} \in \mathcal{C}$ and define a reference measure $\mu(\mathbf{x})$:

$$\mathbf{x} \in \mathcal{C}, \quad \mathbf{m} = \int_{\mathcal{C}} \mathbf{x} \, d\mu(\mathbf{x}), \quad (7)$$

where \mathbf{m} is the mean value of \mathbf{x} under this reference measure.

- Consider \mathbf{x} as the mean value of a random vector \mathbf{X} for which you assume a probability distribution P :

$$\mathbf{x} = \mathbb{E}_P \{ \mathbf{X} \} = \int_{\mathcal{C}} \mathbf{x} \, dP(\mathbf{x}) \quad (8)$$

and the data \mathbf{y} as exact equality constraints on it:

$$\mathbf{y} = \mathbf{Ax} = \mathbf{A} \mathbb{E}_P \{ \mathbf{X} \} = \int_{\mathcal{C}} \mathbf{Ax} \, dP(\mathbf{x}). \quad (9)$$

- Determine the distribution P by:

$$\text{maximize} \quad - \int_{\mathcal{C}} \log \frac{dP(\mathbf{x})}{d\mu(\mathbf{x})} \, dP(\mathbf{x}) \quad \text{s.t.} \quad \mathbf{y} = \mathbf{Ax} = \mathbf{A} \mathbb{E}_P \{ \mathbf{X} \}. \quad (10)$$

The solution is calculated via Lagrangian:

$$\begin{aligned} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) &= \int_{\mathcal{C}} \left[\log \frac{dP(\mathbf{x})}{d\mu(\mathbf{x})} - \sum_{i=1}^M \lambda_i (y_i - [\mathbf{Ax}]_i) \right] dP(\mathbf{x}) \\ &= \int_{\mathcal{C}} \left[\log \frac{dP(\mathbf{x})}{d\mu(\mathbf{x})} - \boldsymbol{\lambda}^t (\mathbf{y} - \mathbf{Ax}) \right] dP(\mathbf{x}) \end{aligned} \quad (11)$$

and is given by:

$$dP(\mathbf{x}, \boldsymbol{\lambda}) = \exp [\boldsymbol{\lambda}^t [\mathbf{Ax}] - \log Z(\boldsymbol{\lambda})] \, d\mu(\mathbf{x}), \quad (12)$$

where

$$Z(\boldsymbol{\lambda}) = \int_{\mathcal{C}} \exp [\boldsymbol{\lambda}^t [\mathbf{Ax}]] \, d\mu(\mathbf{x}). \quad (13)$$

The Lagrange parameters are calculated by searching the unique solution (if exists) of the following system of non linear equations:

$$\frac{\partial \log Z(\boldsymbol{\lambda})}{\partial \lambda_i} = y_i, \quad i = 1, \dots, M. \quad (14)$$

- The solution to the inverse problem is then defined as the expected value of this distribution:

$$\hat{\mathbf{x}}(\boldsymbol{\lambda}) = \int \mathbf{x} \, dP(\mathbf{x}, \boldsymbol{\lambda}). \quad (15)$$

These steps are very formal. In fact, it is possible to determine $\hat{\mathbf{x}}(\boldsymbol{\lambda})$ in a more direct manner. Using the following notations:

$$\mathbf{s} = \mathbf{A}^t \boldsymbol{\lambda}, \quad G^*(\mathbf{s}) = \log Z(\mathbf{s}) = \log \int_{\mathcal{C}} \exp[\mathbf{s}^t \mathbf{x}] \, d\mu(\mathbf{x}), \quad (16)$$

and

$$H(\mathbf{x}) = \max_{\mathbf{s}} \{\mathbf{s}^t \mathbf{x} - G^*(\mathbf{s})\}, \quad D(\boldsymbol{\lambda}) = \boldsymbol{\lambda}^t \mathbf{y} - G^*(\mathbf{A}^t \boldsymbol{\lambda}) \quad (17)$$

it is shown that:

$$\hat{\boldsymbol{\lambda}} = \arg \max_{\boldsymbol{\lambda}} \{D(\boldsymbol{\lambda})\} \quad (\text{Dual criterion}) \quad (18)$$

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x} \in \mathcal{C}} \{H(\mathbf{x})\} \quad \text{s.t.} \quad \mathbf{y} = \mathbf{A}\mathbf{x} \quad (\text{Primal criterion}), \quad (19)$$

$$\hat{\mathbf{x}}(\mathbf{s}) = \frac{dG^*(\mathbf{s})}{d\mathbf{s}} \quad (\text{Explicit relation}), \quad (20)$$

where:

- Functions G and H depend on the reference measure $\mu(\mathbf{x})$;
- $D(\boldsymbol{\lambda})$ is the dual criterion which is function of the data and the function G ;
- $H(\mathbf{x}) = H(\mathbf{x}, \mathbf{m})$ is the primal criterion which is a distance measure between \mathbf{x} and \mathbf{m} which means:
 - $H(\mathbf{x}, \mathbf{m}) \geq 0$, and $H(\mathbf{x}, \mathbf{m}) = 0$ iff $\mathbf{x} = \mathbf{m}$;
 - $H(\mathbf{x}, \mathbf{m})$ is differentiable and convex on \mathcal{C} ;
 - $H(\mathbf{x}, \mathbf{m}) = \infty$ if $\mathbf{x} \notin \mathcal{C}$.

Now, to be able to go a little more in details, let assume that the reference measure is separable:

$$\mu(\mathbf{x}) = \prod_{j=1}^N \mu_j(x_j) \quad (21)$$

then, we have:

$$dP(\mathbf{x}, \boldsymbol{\lambda}) = \prod_{j=1}^N dP_j(x_j, \lambda_j) \quad (22)$$

and

$$G(\mathbf{s}) = \sum_j g_j(s_j), \quad H(\mathbf{x}, \mathbf{m}) = \sum_j h_j(x_j, m_j), \quad \hat{x}_j = g'_j(s_j). \quad (23)$$

Replacing $\mathbf{s} = \mathbf{A}^t \boldsymbol{\lambda}$ we obtain:

$$G(\boldsymbol{\lambda}) = \sum_j g_j([\mathbf{A}^t \boldsymbol{\lambda}]_j), \quad H(\mathbf{x}, \mathbf{m}) = \sum_j h_j(x_j, m_j), \quad \hat{x}_j = g'_j([\mathbf{A}^t \hat{\boldsymbol{\lambda}}]_j), \quad (24)$$

where h_j and g_j depend on the reference measure μ_j :

- g_j is the log Laplace transform (Cramer transform) of μ_j :

$$g(s) = \log \int \exp[sx] d\mu(x);$$

- h_j is the convex conjugate of g_j : $h(x) = \max_{\mathbf{s}} \{sx - g(s)\}$.

Let give some examples:

	$\mu_j(x)$	$g_j(s)$	$h_j(x, m)$
Gaussian:	$\exp\left[-\frac{1}{2}(x-m)^2\right]$	$\frac{1}{2}(s-m)^2$	$\frac{1}{2}(x-m)^2$
Poisson:	$\frac{m^x}{x!} \exp[-m]$	$\exp[m-s]$	$-x \log \frac{x}{m} + m - x$
Gamma:	$x^{\alpha-1} \exp\left[-\frac{x}{m}\right]$	$\log(s-m)$	$-\log \frac{x}{m} + \frac{x}{m} - 1$

When $\mu(\mathbf{x})$ is not separable it is very difficult to do the calculation more in details, excepted the Gaussian case $\mu(\mathbf{x}) = \mathcal{N}(\mathbf{m}, \mathbf{R}_x)$, where we have:

$$H(\mathbf{x}, \mathbf{m}) = -\frac{1}{2}(\mathbf{x} - \mathbf{m})^t \mathbf{R}_x^{-1}(\mathbf{x} - \mathbf{m}), \quad G(\boldsymbol{\lambda}) = -\frac{1}{2}\|\boldsymbol{\lambda}\|^2, \quad D(\boldsymbol{\lambda}) = \boldsymbol{\lambda}^t \mathbf{y} + \frac{1}{2}\|\boldsymbol{\lambda}\|^2. \quad (25)$$

(See however a new presentation of the method in [?] trying to extend the method for taking account of the correlations.)

2.2. EXTENSIONS

How to account for the noise: Two approaches have been developed in [?, ?, ?, ?]:

- Replacing the exact equality constraints $\mathbf{y} = \mathbf{A}\mathbf{x}$ by the following inequalities:

$$|y_i - [\mathbf{A}\mathbf{x}]_i| < \epsilon, \quad \text{or} \quad \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2 < \sigma^2 \quad (26)$$

and using the duality relations they showed:

$$\begin{cases} \hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \{H(\mathbf{x})\} \text{ s.t. } |y_i - [\mathbf{A}\mathbf{x}]_i| < \epsilon, \text{ with } H(\mathbf{x}) = \sum_j h_j(x_j) \\ \hat{\boldsymbol{\lambda}} = \arg \max_{\boldsymbol{\lambda}} \{\tilde{D}(\boldsymbol{\lambda})\} \text{ with } \tilde{D}(\boldsymbol{\lambda}) = D(\boldsymbol{\lambda}) + \alpha\|\boldsymbol{\lambda}\|^2 \end{cases} \quad (27)$$

where α depends on ϵ or on σ^2 .

- Replacing $\mathbf{y} = \mathbf{A}\mathbf{x}$ by $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{n}$ and rewriting it as follows:

$$\mathbf{y} = [\mathbf{A}|\mathbf{I}] \begin{bmatrix} \mathbf{x} \\ \mathbf{n} \end{bmatrix} \longrightarrow \mathbf{y} = \tilde{\mathbf{A}} \tilde{\mathbf{x}} \quad (28)$$

and assuming $\mu(\tilde{\mathbf{x}}) = \mu_x(\mathbf{x})\mu_n(\mathbf{n})$ they showed:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x} \in \mathcal{C}} \{ \mathcal{Q}(\mathbf{y} - \mathbf{A}\mathbf{x}) + \alpha H(\mathbf{x}) \} \quad (29)$$

$$\text{with} \quad H(\mathbf{x}) = \sum_{j=1}^N h_j(x_j), \quad \text{and} \quad \mathcal{Q}(\mathbf{z}) = \sum_{i=1}^M q_i(z_i). \quad (30)$$

Here also $h_j(x_j)$ and $q_i(z_i)$ depend on the reference measures $\mu_x(\mathbf{x})$ and $\mu_n(\mathbf{x})$. The determination of α is not discussed.

3. Bayesian approach

3.1. BASICS

The different steps of this approach are now well-known:

- From the observation model and the hypothesis (prior knowledge) on the noise derive the likelihood $p(\mathbf{y}|\mathbf{x};\boldsymbol{\beta})$;
- From the hypothesis (prior knowledge) on \mathbf{x} derive the prior law $p(\mathbf{x}|\boldsymbol{\theta})$;
- Apply the Bayes rule to obtain $p(\mathbf{x}|\mathbf{y};\boldsymbol{\beta},\boldsymbol{\theta}) = p(\mathbf{y}|\mathbf{x};\boldsymbol{\beta})p(\mathbf{x}|\boldsymbol{\theta})/p(\mathbf{y};\boldsymbol{\beta},\boldsymbol{\theta})$;
- Define an estimation rule via a cost function $c(\mathbf{x},\hat{\mathbf{x}})$ by:

$$\hat{\mathbf{x}}(\mathbf{y};\boldsymbol{\beta},\boldsymbol{\theta}) = \arg \min_{\mathbf{z}} \left\{ \int c(\mathbf{x},\mathbf{z})p(\mathbf{x}|\mathbf{y};\boldsymbol{\beta},\boldsymbol{\theta}) \mathrm{d}\mathbf{x} \right\}. \quad (31)$$

Different cost functions give different estimators:

- Maximum a posteriori (MAP):

$$C(\mathbf{x},\hat{\mathbf{x}}) = 1 - \delta(\mathbf{x} - \hat{\mathbf{x}}) \longrightarrow \hat{\mathbf{x}} = \arg \max_{\mathbf{x}} \{p(\mathbf{x}|\mathbf{y};\boldsymbol{\theta},\boldsymbol{\beta})\}. \quad (32)$$

- Posterior mean (PM):

$$C(\mathbf{x},\hat{\mathbf{x}}) = [\mathbf{x} - \hat{\mathbf{x}}]^t \mathbf{Q} [\mathbf{x} - \hat{\mathbf{x}}]^t \longrightarrow \hat{\mathbf{x}} = \mathbb{E}_{\mathbf{x}|\mathbf{y}} \{ \mathbf{X} \} = \int \mathbf{x} p(\mathbf{x}|\mathbf{y};\boldsymbol{\theta},\boldsymbol{\beta}) \mathrm{d}\mathbf{x}. \quad (33)$$

- Maximum of the Marginal a posteriori (MMAP):

$$C(\mathbf{x},\hat{\mathbf{x}}) = \prod_j 1 - \delta(x_j - \hat{x}_j) \longrightarrow \hat{x}_j = \arg \max_{x_j} \{p(x_j|\mathbf{y};\boldsymbol{\theta})\}, \quad (34)$$

where

$$p(x_j|\mathbf{y};\boldsymbol{\theta}) = \int p(\mathbf{x}|\mathbf{y};\boldsymbol{\theta}) \mathrm{d}x_1 \cdots \mathrm{d}x_{j-1} \cdots \mathrm{d}x_{j+1} \cdots \mathrm{d}x_n. \quad (35)$$

To illustrate this, let consider the case of linear inverse problems $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{n}$ with the following hypothesis:

- \mathbf{n} is zero-mean, white and Gaussian: $\mathbf{n} \sim \mathcal{N}(\mathbf{0}, (1/\beta)\mathbf{I})$ which leads to:

$$p(\mathbf{y}|\mathbf{x}, \beta) \propto \exp \left[-\frac{1}{2}\beta \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2 \right]. \quad (36)$$

- \mathbf{x} is Gaussian: $\mathbf{x} \sim \mathcal{N}(\mathbf{x}_0, (1/\theta)\mathbf{P}_0)$:

$$p(\mathbf{x}|\theta) \propto \exp \left[-\frac{1}{2}\theta [\mathbf{x} - \mathbf{x}_0]^t \mathbf{P}_0^{-1} [\mathbf{x} - \mathbf{x}_0] \right]. \quad (37)$$

Then, using the Bayes rule it is easy to show that

$$\mathbf{x}|\mathbf{y} \sim \mathcal{N}(\hat{\mathbf{x}}, \mathbf{P}) \quad \text{with} \quad \hat{\mathbf{x}} = \mathbf{P}\mathbf{A}^t(\mathbf{y} - \mathbf{A}\mathbf{x}_0), \quad \mathbf{P} = (\mathbf{A}^t\mathbf{A} + \lambda\mathbf{P}_0^{-1})^{-1}. \quad (38)$$

The MAP solution is:

$$\hat{\mathbf{x}} = \arg \max_{\mathbf{x}} \{p(\mathbf{x}|\mathbf{y})\} = \arg \min_{\mathbf{x}} \{J(\mathbf{x})\}, \quad \text{with} \quad J(\mathbf{x}) = Q(\mathbf{x}) + \lambda\phi(\mathbf{x}), \quad (39)$$

where

$$Q(\mathbf{x}) = \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2, \quad \phi(\mathbf{x}) = \mathbf{x}^t \mathbf{P}_0^{-1} \mathbf{x} = \|\mathbf{D}\mathbf{x}\|, \quad \lambda = \frac{\theta}{\beta} \quad (40)$$

Now, relaxing the second hypothesis, i.e; choosing other prior laws for \mathbf{x} we obtain other MAP criteria. Let just note some special interesting cases:

- A Generalized Gaussian law for \mathbf{x} :

$$p(x_j) \propto \exp [-(x_j - m_j)^\alpha]. \quad (41)$$

The related MAP criterion becomes:

$$J(\mathbf{x}) = Q(\mathbf{x}) + \phi(\mathbf{x}) \quad \text{with} \quad \phi(\mathbf{x}) = \sum_j (x_j - m_j)^\alpha. \quad (42)$$

- A Gamma law for \mathbf{x} :

$$x_j \sim \mathcal{G}(\alpha, m_j) \longrightarrow p(x_j) \propto (x_j/m_j)^{-\alpha} \exp [-x_j/m_j]. \quad (43)$$

The related MAP criterion becomes:

$$J(\mathbf{x}) = Q(\mathbf{x}) + \phi(\mathbf{x}) \quad \text{with} \quad \phi(\mathbf{x}) = \sum_j \alpha \log \frac{x_j}{m_j} + \frac{x_j}{m_j}. \quad (44)$$

- A Beta law for \mathbf{x} :

$$x_j \sim \mathcal{B}(\alpha, \beta) \longrightarrow p(x_j) \propto x_j^\alpha (1 - x_j)^\beta. \quad (45)$$

The related MAP criterion becomes:

$$J(\mathbf{x}) = Q(\mathbf{x}) + \phi(\mathbf{x}) \quad \text{with} \quad \phi(\mathbf{x}) = \alpha \sum_j \log x_j + \beta \sum_j \log(1 - x_j). \quad (46)$$

- A Poisson law for \mathbf{x} :

$$p(x_j) \propto \frac{m_j^{x_j}}{x_j!} \exp[-m_j]. \quad (47)$$

The related MAP criterion becomes:

$$J(\mathbf{x}) = Q(\mathbf{x}) + \phi(\mathbf{x}) \quad \text{with} \quad \phi(\mathbf{x}) = - \sum_j x_j \log \frac{x_j}{m_j} + (x_j - m_j). \quad (48)$$

- Markovian models for \mathbf{x} :

$$J(\mathbf{x}) = Q(\mathbf{x}) + \phi(\mathbf{x}) \quad \text{with} \quad \phi(\mathbf{x}) = \alpha \sum_j \sum_{i \in N_j} V(x_j, x_i). \quad (49)$$

3.2. EXTENSIONS

The Bayesian approach can be exactly applied when all the direct (prior) probability laws $(p(\mathbf{y}|\mathbf{x}, \boldsymbol{\beta}))$ and $p(\mathbf{x}|\boldsymbol{\theta})$ are assigned. Even, choosing an appropriate law is done in general by hand, another difficulty is to determine their parameters $(\boldsymbol{\beta}, \boldsymbol{\theta})$. This problem has been addressed by many authors and the subject is an active area in statistics. See [?, ?, ?, ?], [?, ?, ?, ?, ?] and also [?, ?, ?].

All these methods can mainly be divided in three main families:

- Generalized MAP: In this approach one tries to estimate both the hyperparameters and the unknown variables \mathbf{x} directly from the data by defining:

$$(\hat{\mathbf{x}}, \hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\beta}}) = \arg \max_{(\mathbf{x}, \boldsymbol{\theta}, \boldsymbol{\beta})} \{p(\mathbf{x}, \boldsymbol{\theta}, \boldsymbol{\beta}|\mathbf{y})\} \quad (50)$$

where

$$p(\mathbf{x}, \boldsymbol{\theta}, \boldsymbol{\beta}|\mathbf{y}) \propto p(\mathbf{y}|\mathbf{x}, \boldsymbol{\beta}) p(\mathbf{x}|\boldsymbol{\theta}) p(\boldsymbol{\theta}) p(\boldsymbol{\beta}) \quad (51)$$

and where $p(\boldsymbol{\theta})$ and $p(\boldsymbol{\beta})$ are appropriate prior laws. Many authors used the non informative prior law for them.

- Marginalization: In this approach one tries to estimate first the hyperparameters by marginalizing over the unknown variables \mathbf{x} :

$$p(\boldsymbol{\theta}, \boldsymbol{\beta}|\mathbf{y}) \propto p(\boldsymbol{\beta}) p(\boldsymbol{\theta}) \int p(\mathbf{y}|\mathbf{x}, \boldsymbol{\beta}) p(\mathbf{x}|\boldsymbol{\theta}) d\mathbf{x} \quad (52)$$

and then, using them in the estimation of the unknown variables \mathbf{x} :

$$(\hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\beta}}) = \arg \max_{(\boldsymbol{\theta}, \boldsymbol{\beta})} \{p(\boldsymbol{\theta}, \boldsymbol{\beta}|\mathbf{y})\} \quad (53)$$

- Nuisance parameters: In this approach the hyperparameters are considered as the nuisance parameters, so, marginalized:

$$p(\mathbf{x}|\mathbf{y}) = \int p(\mathbf{y}, \mathbf{x}, \boldsymbol{\theta}, \boldsymbol{\beta}) d\boldsymbol{\theta} d\boldsymbol{\beta} \quad (54)$$

and the unknown variables \mathbf{x} are estimated by:

$$\hat{\mathbf{x}} = \arg \max_{\mathbf{x}} \{p(\mathbf{x}|\mathbf{y})\} \quad (55)$$

To see some more discussions and different possible implementations of these approaches see [?].

4. Discussed points in each approach

4.1. MEM:

- Choice of \mathcal{C} and $\mu(\mathbf{x})$:
 - \mathcal{C} must be a convex set, such as: \mathbf{R}^N , \mathbf{R}_+^N , $[a, b]^N$
 - Up to now, the whole analysis can be done for separable measures.
 - The only reference measures μ_j which permits to go through all the steps are those for which we have analytical expression for the Laplace transform of their logarithms.
- Accounting for the noise:

In the first approach only the support and the energy of the noise is used. In the second approach we have more choices via the reference measures $\mu_n(\mathbf{n})$, but determination of α stays adhoc. In fact, in general, when the reference measures $\mu_n(\mathbf{n})$ and $\mu_x(\mathbf{x})$ depend on any parameters, this approach lacks any tool to determine them.
- Effective calculation of the solution:

No problem, and more, this is probably the main interest of the approach which defines the solution, by construction, as the optimizer of a convex criterion.
- Characterization of the solution:

A sensitivity analysis has been proposed by [?], but, in my opinion, this is not enough to characterize a solution to an inverse problem.

It is not easy to use the notions of variance or covariance of the solution, because this approach does not define a posterior distribution for the solution.
- Possibility of the extension of the approach:

I have not yet seen any extension of this approach to non linear inverse problems, or linear inverse problems in which the operator depend on unknown parameters, such as blind deconvolution or antenna array processing;

The fact that we have to choose a convex set \mathcal{C} on which the solution is defined excludes the inverse problems in which we know a priori that the solution is discrete-valued (binary, *n*-ary images for example). This excludes the use of this approach in image segmentation or communication inverse problems (canal equalization, blind deconvolution, etc.).

4.2. BAYES:

- Choice of $p(\mathbf{x}|\boldsymbol{\beta})$:

$p(\mathbf{x}|\boldsymbol{\beta})$ can be chosen separable or not. Evidently, separable $p(\mathbf{x}|\boldsymbol{\beta})$ (Entropic prior laws) simplifies the calculations. Accounting for correlations is easily done via Markovian models; In both cases (Entropic or Markovian prior laws), there are some tools for choosing them either by physical considerations, or by scale invariance arguments [?, ?, ?, ?, ?].
- Choice of the cost function or equivalently of an estimator MAP, PM, MMAP:

This choice is done more on the basis of cost calculation. MAP calculation needs, in general, global optimization, but does not need any integration. MP or MMAP needs multidimensional integration, so in general, greater cost. However, there are approximate calculation techniques based on Monte Carlo

methods and Gibbs sampling.

- Effective calculation of different solutions:

For MAP estimate, when the posterior law is unimodal, we can use any gradient descent based method, but if this is not the case, there are two categories of methods: Simulated Annealing or Deterministic relaxation (GNC). For more discussions on Bayesian calculations see [?] in this volume.

5. Comparisons and discussions

The following main items are discussed:

- In MEM, the unknowns \mathbf{x} are considered as the mean values of a random vector \mathbf{X} for which a prior probability measure $d\mu(\mathbf{x})$ is defined.
- In BAYES, the unknowns \mathbf{x} are considered as a sample of a random vector \mathbf{X} for which a prior probability measure $p(\mathbf{x})$ is defined.
- In MEM, a probability distribution $p(\mathbf{x})$ is defined as the minimizer of the Kullback distance $K(p, \mu)$ subject to the data constraints, and the solution is defined to be $E_p(\mathbf{X})$. What is interesting here is that this solution can equivalently be obtained as the minimizer of a convex criterion $J(\mathbf{x})$ subject to the data constraints, and what is more attractive is that, thanks to the convex analysis, this solution can also be obtained as the stationary point of a dual criterion which can easily be calculated numerically.
- In BAYES, the posterior law $p(\mathbf{x}|\mathbf{y})$ is calculated using the Bayes' rule. In fact, the data \mathbf{y} are considered as a sample of a random vector \mathbf{Y} for which we can define a conditional probability law $p(\mathbf{y}|\mathbf{x})$ which, when used in conjunction with the prior $p(\mathbf{x})$ in the Bayes' rule will give us the posterior law, from which we can define an estimator. One of these estimators is the posterior mean $E_p(\mathbf{X})$, but others can also be defined. This posterior law is used not only to define an estimate (a solution), but also to calculate any other probabilistic information about that solution.
- In MEM, in its original version, the data are considered as the exact linear constraints. The uncertainty on the data are not considered, and consequently, the uncertainty on the solution is not handled. However, some extensions are recently presented to take account of the errors on the data and to calculate the sensitivity of the solution to these errors.
- In BAYES, the errors are naturally considered through $p(\mathbf{y}|\mathbf{x})$ and the uncertainty of the solution through the posterior probability $p(\mathbf{x}|\mathbf{y})$. Naturally then, we can compare the information content of the data and the prior model using their entropies. We can also measure the relative information content of the posterior to prior model by $K(p(\mathbf{x}|\mathbf{y}), p(\mathbf{x}))$.
- In MEM, even in their extended versions, it is not easy to handle with the hyperparameters. In BAYES, there are the necessary tools to handle them.
- In MEM, one can not yet handle with non linear problems. This is not the case of the BAYES.

As a final conclusion, we have to mention that, even if the two approaches are different, they can, in some cases result to the same definition of the solution as

the minimizer of the same criterion, and consequently, to give exactly the same numerical solutions to a given inverse problem. However, we can give different interpretations to the obtained result depending on the approach used to reach it. The main objective of this paper was to give a succinct presentation of the two approaches for the resolution of the inverse problems.

It is important to note that the two approaches give different views and interpretations which can be used advantageously for any application. Also, even the Bayesian approach is now really mature, the MEM approach is more recent. So, many of the conclusions I made today may be altered in future. In particular, new presentation of the method by Heinrich et al [?] in this volume will probably give new possibilities to the MEM method and will push the limits of the method to greatest generality.